

**SEPARATE READ AND WRITE SERVERS
IN A DISTRIBUTED FILE SYSTEM**

Inventor(s)

Christos Karamanolis
1657 Belleville Way, Apt. J
Sunnyvale, CA 94087

Daniel A. Muntz
10301 Vicksburg Drive
Cupertino, CA 95014

Mallik Mahalingam
620 Park View Drive, #105
Santa Clara, CA 95054

Zheng Zhang
3520 Casabella Court
San Jose, CA 95148

Assignee

Hewlett Packard Company

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

SEPARATE READ AND WRITE SERVERS IN A DISTRIBUTED FILE SYSTEM

FIELD OF THE INVENTION

The present invention generally relates to distributed file systems, and more particularly to separate servers for reading and writing data in a distributed file system.

BACKGROUND

Distributed file systems are generally employed for storage of large quantities of data and to reduce input/output (I/O) bottlenecks where there are many requests made for file access. In a distributed file system, the file data is spread across multiple data processing systems. File system control and management of file system meta-data is distributed in varying degrees in different systems.

A desirable characteristic of many distributed file systems is scalability. Scalability is a characteristic that refers to the ease with which a distributed file system can be expanded to accommodate increased data access needs or increased storage needs. For example, as additional users are granted access to the distributed file system, new storage servers may be introduced, and the requests of the additional users may be further spread across the old servers and new servers. The scalability of any distributed file system is limited or enhanced by the system design.

Caching is a feature that is commonly used to reduce data access times and to enhance the scalability of distributed file systems. However, caching requires additional management software to address data locking and data consistency issues. Thus, caching introduces additional complexity and overhead into a distributed file system.

Another approach that addresses scalability is the provision of dual paths for access to file data and access to meta-data. In this approach, the meta-data is managed on a server that is separate from the storage servers. However, this approach may create a bottleneck at the meta-data server and thereby restrict scalability.

A system and method that address the aforementioned problems, as well as other related problems, are therefore desirable.

SUMMARY OF THE INVENTION

2 In various embodiments, a system and method are provided for implementing a
3 distributed file system in which read requests are processed by dedicated read servers
4 and write requests are processed by a dedicated write server. In various embodiments,
5 read requests are separated from write requests and processed by dedicated read
6 servers. A plurality of read servers are coupled to the client applications and each read
7 server reads file data from the distributed file system and returns the file data to the
8 client applications. A write server writes data to the distributed file system. Various
9 embodiments are described for separating read requests from write requests and
10 transmitting read requests to the read servers and write requests to the write server.

11 It will be appreciated that various other embodiments are set forth in the
12 Detailed Description and Claims which follow.

BRIEF DESCRIPTION OF THE DRAWINGS

15 Various aspects and advantages of the invention will become apparent upon
16 review of the following detailed description and upon reference to the drawings in
17 which:

18 FIG. 1 is a functional block diagram that illustrates the flow of file access
19 requests and file data in accordance with one embodiment of the invention;

20 FIG. 2 is a functional block diagram of a system where separate read and write
21 servers provide access to file data;

22 FIG. 3 is a functional block diagram of a system where separate read and write
23 servers provide access to file data and a dedicated load balancer distributes the file
24 access requests;

25 FIG. 4 is a functional block diagram of software components hosted by a read
26 server;

27 FIG. 5 is a functional block diagram of software components hosted by a write
28 server;

29 FIG. 6 is a flowchart that illustrates a process performed by the load balancer in
30 processing file access requests;

31 FIG. 7 is a flowchart of a process performed by a read server in processing
32 requests from a distributed file system client interface;

33 FIG. 8 is a flowchart of a process performed by the write server in writing data
34 to a file; and

1 FIGs. 9A, 9B, and 9C, which illustrate successive states of file storage in
2 processing a write request.

3

4 **DETAILED DESCRIPTION**

5 Various embodiments of the present invention are described in terms of specific
6 functions implemented on specific data processing systems. Those skilled in the art
7 will appreciate, however, that various alternative arrangements of data processing
8 systems and various alternative data processing system architectures could be used to
9 implement the invention.

10 FIG. 1 is a functional block diagram that illustrates the flow of file access
11 requests and file data in accordance with one embodiment of the invention. System
12 100 includes a plurality of clients 102-104, a plurality of read servers 106-108, a write
13 server 110, and a distributed file system 112. The clients 102-104 are data processing
14 systems that host client applications (not shown) that issue requests to read data from
15 and write data to storage that is managed by distributed file system 112. Read requests
16 refer to file access operations where file system meta-data, file meta-data, or file data
17 are read from storage, and write requests refer to file access operations where file
18 system meta-data, file meta-data, or file data are written to storage.

19 To improve scalability and performance, read requests and write requests are
20 processed by different servers. The clients 102-104 send read requests to the read
21 servers 106-108 for processing, and write requests are sent to write server 110. In
22 applications where read activity is much greater than write activity, the separation of
23 read and write servers supports scalability to service more read requests. Measurements
24 from commercial systems indicate that read operations are typically more than 90% of
25 the total operations to a distributed file system. Thus, additional read servers can be
26 coupled to the client applications and to the distributed file system 112 to handle more
27 read requesters. The addition of read servers does not require any reconfiguration of
28 the distributed file system and can be transparent to the user application.

29 In one embodiment, the particular read server to which a client application
30 sends a read request is selected in a manner that balances the processing load between
31 the read servers. Each read server provides access to all the addressable file storage for
32 each of the coupled client applications.

33 In an example embodiment, the read servers 106 and 108 and write server 110
34 are implemented as conventional network file system (NFS) servers that are coupled to

1 a conventional distributed file system 112 and hosted on separate data processing
2 systems. In another embodiment, the read servers are adapted to receive all file access
3 requests and forward write requests (not shown) to the write server 110. Those skilled
4 in the art will recognize that various alternative remote and distributed file systems
5 could be adapted to operate in accordance with the present invention.

6 By separating the read and write requests, the system 100 is scalable to process
7 read requests. Since the read servers do not write any data to storage they do not
8 require consistency control of the data, and additional read servers can be added
9 without incurring extra overhead on the other read servers. Since there is only one
10 write server, the overhead associated with maintaining data consistency between
11 multiple writers is eliminated.

12

13 FIG. 2 is a functional block diagram of a system 150 where separate read and
14 write servers provide access to file data. Client systems 102 and 104, read servers 106
15 and 108, and write server 110 are coupled to network 152. Network 152 is a
16 conventional data communications network through which the servers and clients
17 interact. Storage area network 154 includes the physical storage media on which the
18 file data are stored. Storage Area Networks (SANs) include three types of hardware
19 components: fiber channel switches (e.g., BROCADE FC-SAN Switches), fiber
20 channel adaptors for hosts (e.g., Qlogic FC-SAN host adaptors), and special disks or
21 disk arrays (e.g., Hewlett-Packard's XP512 Disk Arrays). Special software is required
22 to manage and configure SANs, e.g., McData's SAN Management and IBM's Tivoli.

23 Each of client systems 102 and 104 hosts a client application and an interface to
24 the distributed file system. For example, client 102 hosts client application 156 and
25 distributed file system (DFS) client interface 158. Other than file access requests made
26 by the client application 156, the application-specific functions of the client application
27 are beyond the scope of the present invention.

28 The DFS client interface 158 is implemented with functional extensions to
29 conventional DFS client software. For example, in one embodiment NFS-client
30 software is extended with functions that separate read requests from write requests and
31 send the requests to the appropriate servers. In another embodiment, the DFS client
32 interface 158 is implemented with conventional NFS-client software, and the read
33 servers 106 and 108 are adapted to separate read requests from write requests. In the
34 latter embodiment, the read servers forward write requests to the write server 110.

1 Patent/ application number *****, entitled, "EXTENDING A STANDARD-BASED
2 REMOTE FILE ACCESS PROTOCOL AND MAINTAINING COMPATIBILITY
3 WITH A STANDARD PROTOCOL STACK" by Karamanolis et al., filed on January
4 31, 2001, and assigned to the assignee of the present invention, describes yet another
5 embodiment for implementing the DFS client interface and is hereby incorporated by
6 reference. It will be appreciated that other standards-based or proprietary distributed
7 file systems can be adapted in accordance with the teachings of present invention.

8 In another embodiment of the invention, the DFS client interface 158 includes
9 functionality that distributes read requests between the read servers 106 and 108 in
10 order to balance the processing load between the read servers. For example, a round-
11 robin or other well known load distribution function can be used to balance read
12 requests between the read servers.

13

14 FIG. 3 is a functional block diagram of a system 160 where separate read and
15 write servers provide access to file data and a dedicated load balancer distributes the
16 file access requests. DFS client interface 162 sends read and write requests to load
17 balancer 164 instead of addressing the read and write servers directly. In other
18 respects, DFS interface 162 is implemented as described above.

19 In one embodiment, load balancer 164 is implemented with a conventional
20 content switch that is coupled to network 152. The load balancer 164 is an application
21 layer switch (i.e., layer 7). Application layer switches that currently support switching
22 for URLs can be programmed or configured to support distributed file system access.
23 The load balancer 164 is configured to receive read and write requests from the DFS
24 client interface 162 components on each of the clients 102-104. In a first embodiment,
25 the load balancer distributes read and write requests to the read servers, and the read
26 servers are configured to forward the write requests to the write server 110. In another
27 embodiment, the load balancer distributes read requests to the read servers and
28 forwards write requests to the write server. Based on the function code present in a file
29 access request, the load balancer distinguishes between read and write requests.

30 Load balancer 164 attempts to evenly distribute the processing load associated
31 with servicing read requests between the read servers. In one embodiment, a round-
32 robin method is used to distribute the requests to the read servers. More sophisticated
33 approaches may be employed in other embodiments. For example, the load balancer
34 can examine each read request for the quantity of data requested and use the

1 combination of the quantity of data and number of outstanding read requests to evenly
2 distribute the workload. In another embodiment, each read server reports its workload
3 to the load balancer, and the load balancer uses the relative current workloads of the
4 read servers in distributing read requests.

5

6 FIG. 4 is a functional block diagram of software components hosted by a read
7 server 182. Read server 182 is a conventional data processing system having
8 computational and input/output capacity that depend on application requirements. In
9 various embodiments, the software components are conventional or have extensions
10 added to conventional software components.

11 DFS server 184 receives file access requests from the client application 156. In
12 one embodiment, the DFS server is implemented with conventional server software for
13 a distributed file system, for example, NFS server software. If the DFS client interface
14 158 or load balancer 164 sends only read requests to the DFS server, the DFS server
15 processes only read requests and commercially available DFS server software is used.
16 In another embodiment, DFS client interface 158 or load balancer 164 sends both read
17 and write requests to the DFS server, and the DFS server is configured to forward write
18 requests to the write server 110.

19 Physical file system 186 is also implemented with conventional software. For
20 example, the physical file system can be implemented with the Ext2 system of Linux or
21 the FFS of BSD Unix. Alternatively, proprietary software such as NTFS from
22 Microsoft, XFS from Silicon Graphics, or WAFL from Network Appliances, may be
23 used to implement the physical file system..

24

25 FIG. 5 is a functional block diagram of software components hosted by a write
26 server. Write server 192 is a conventional data processing system having computational
27 and input/output capacity that depend on application requirements. Comparable to the
28 read server 182, the software components are based on conventional software
29 components.

30 DFS server 184 processes write requests from the client application 156. The
31 DFS server 184 is adapted to interface with data consistency control element 194.
32 Since the read and write servers have access to the same virtual storage, when file data
33 and meta-data are modified the write server must ensure that the data are modified in a
34 consistent manner. That is, the data and meta-data read by the read servers must be

1 consistent. "Meta-data" refers to information that describes the file system and
2 information that describes each file. For example, meta-data includes status
3 information, permission levels, physical storage location descriptions, symbolic names,
4 etc.

5 The data consistency control logic 194 assumes that the client application 156
6 does not immediately require the most recent data. Once the data consistency control
7 194 has stored the new meta-data and file data in a consistent state, the new data is
8 accessible to the read servers.

9 As described below in figures 8 and 9A – 9C, the write server imposes a strict
10 order of operations in accessing the physical storage (e.g., disk) when servicing a write
11 request. This requires support from the physical file system because the physical file
12 system controls data consistency. In one embodiment, the physical file system provides
13 the interface and mechanisms to specify such order requirements. Alternatively,
14 extensions to the physical file system, for example, data consistency control 194,
15 control the order of operations.

16

17 FIG. 6 is a flowchart that illustrates a process performed by the load balancer
18 164 in processing file access requests. At step 302, a file access request is received via
19 network 152 from a DFS client interface 162. In one embodiment, the load balancer is
20 configured to process only read requests (e.g., the DFS client interface separates read
21 requests from write requests), and in another embodiment, the load balancer is
22 configured to process both read and write requests (e.g., the DFS client interface
23 forwards both read and write requests to the load balancer).

24 In the embodiment where the load balancer receives only read requests, the
25 process continues at step 304 where a read server is selected. As described above, the
26 load balancer attempts to balance the workload between the read servers. For example,
27 the load balancer implements a round-robin or other known load balancing algorithm.
28 At step 306, the request is forwarded to the selected read server, and control returns to
29 step 302 to process the next request.

30 In the embodiment where the load balancer receives both read and write
31 requests, the process is directed from step 302 to step decision step 308. At decision
32 step 308, the load balancer checks whether the request is a read request or a write
33 request. For read requests, the process is directed to step 304 and the read request is
34 processed as described above. For write requests, the process is directed to step 310

1 where the write request is sent to the write server 110. The process then returns to step
2 302 to process the next request.

3

4 FIG. 7 is a flowchart of a process performed by a read server in processing
5 requests from a DFS client interface. At step 352, a file access request is received via
6 network 152 from a DFS client interface 158. In one embodiment, the read server is
7 configured to process only read requests (e.g., where the DFS client interface separates
8 read requests from write requests), and in another embodiment, the read server is
9 configured to process read requests and forward write requests to the write server (e.g.,
10 where the DFS client interface forwards both read and write requests to the read
11 servers).

12 In the embodiment where the read server receives only read requests, the
13 process continues at step 354 where a read server is selected as described above. At
14 step 356, the request is forwarded to the selected read server, and control returns to step
15 302 to process the next request.

16 In the embodiment where the read server receives both read and write requests,
17 the process is directed from step 352 to step decision step 358. At decision step 358,
18 the read server checks whether the request is a read request or a write request. For read
19 requests, the process is directed to step 354 and the read request is processed as
20 described above. For write requests, the process is directed to step 360 where the write
21 request is sent to the write server 110. The process then returns to step 352 to process
22 the next request.

23

24 FIG. 8 is a flowchart of a process performed by the write server in writing data
25 to a file. The process of FIG. 8 is described in conjunction with the block diagrams of
26 FIGs. 9A, 9B, and 9C, which illustrate successive states of file storage in processing a
27 write request. At step 402, a write request is received, either from a read server, a load
28 balancer, or from a DFS client interface, depending on the implementation. The i-node
29 for the file referenced in the write request is read at step 404.

30 FIG. 9A illustrates the initial state of file storage 500. File storage 500 includes
31 a block bitmap 502, an i-node 504, multiple indirect blocks 506-508, and multiple data
32 blocks 510-512. The i-node includes pointers to the indirect blocks and pointers to a
33 number of data blocks, for example data block 510. The indirect blocks include
34 pointers to data blocks, for example, indirect block 508 references data block 512.

1 While not shown, it will be appreciated that the file system also includes double and
2 triple indirect blocks as understood by those skilled in the art.

3 The block bitmap 502 indicates which blocks of file storage 500 have file data
4 and meta-data stored therein, and which blocks are available. The i-node 504 contains
5 information that describes the file, for example, a symbolic name, timestamps, and
6 access permissions.

7 The information from i-node 504 that is read into memory of the write server at
8 step 404 is shown as block 514. Step 406 conditionally reads blocks of file data if the
9 write request involves updating presently allocated indirect and data blocks.

10 At step 408, the file data in the write request is used to update data blocks in the
11 memory of the write server. In addition, the i-node and block bitmap are updated in the
12 memory of the write server if necessary (FIG. 9B, 514').

13 At step 412, the file data from the memory of the write server is written to
14 newly allocated data blocks in file storage 500. For example, in FIG. 9B indirect block
15 508 and data block 512 are updated and written to file storage 500 as indirect block
16 508' and data block 512'. In addition, FIG. 9B illustrates a new data block 516 that is
17 written to the file storage. Note that the current i-node 504 still references the old
18 indirect block 508, which references old data block 512. Thus, the read servers
19 continue to have a consistent view of the file data while the write server is processing
20 write requests, even though the data may not be current.

21 At step 414, the portion of file storage 500 having the block bitmap 502 and i-
22 node 504 are locked, and the updated I-node 504' and block bitmap 502' (FIG. 9C) are
23 written to file storage 500 at step 416. Any old data or indirect blocks are freed at step
24 416. The block bitmap and i-node areas are unlocked at step 418. FIG. 9C illustrates
25 the state of file storage 500 after the block bitmap and i-node have been updated. The
26 updated i-node 504' references indirect block 508'. Thus, the read servers have access
27 to the new data after the write server completes the i-node update.

28 At step 420, a response is returned to the DFS client interface, and the process
29 returns to step 402 to receive the next write request.

30 The present invention is believed to be applicable to a variety of distributed and
31 remote files systems and has been found to be particularly applicable and beneficial
32 with NFS-type file systems. Those skilled in the art will appreciate that the invention is
33 not limited to NFS-type file systems, and other aspects and embodiments of the present
34 invention will be apparent from consideration of the specification and practice of the

1 invention disclosed herein. It is intended that the specification and illustrated
2 embodiments be considered as examples only, with a true scope and spirit of the
3 invention being indicated by the following claims.

4

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100